

Algorithmic Composition of Melodies with Deep Recurrent Neural Networks

Florian Colombo, Samuel P. Muscinelli, Alexander Seeholzer, Johanni Brea
and Wulfram Gerstner

Laboratory of Computational Neurosciences. Brain Mind Institute. École
Polytechnique Fédérale de Lausanne
`florian.colombo@epfl.ch`

Abstract. A big challenge in algorithmic composition is to devise a model that is both easily trainable and able to reproduce the long-range temporal dependencies typical of music. Here we investigate how artificial neural networks can be trained on a large corpus of melodies and turned into automated music composers able to generate new melodies coherent with the style they have been trained on. We employ gated recurrent unit networks that have been shown to be particularly efficient in learning complex sequential activations with arbitrary long time lags. Our model processes rhythm and melody in parallel while modeling the relation between these two features. Using such an approach, we were able to generate interesting complete melodies or suggest possible continuations of a melody fragment that is coherent with the characteristics of the fragment itself.

Keywords: algorithmic composition, generative model of music, machine learning, deep recurrent neural networks

1 Introduction

The algorithmic formalization of musical creativity and composition, foreseen as early as the 19th century [1], has come to fruition in the recent decades with the advent of modern computer algorithms [2].

Formally, a melody can be seen as a sample from a potentially very sophisticated probability distribution over sequences of notes [2–5]. For monophonic music, probability distributions could be given by Markov chains, where the probability of the next note depends only on the current note and the k last notes [4]. Markov chain models, however, do not capture the long-range temporal structure inherent in music. For example, even a simple melody such as *Brother John* is structured in four patterns, each repeated twice, with the first and last ones starting with the same notes (see Fig. 1). Taking only the last few notes into account is thus not enough to produce the sequence – rather does the progression on the long timescale of bars dictate the sequences of notes.

Such rich temporal structure can be captured by models that rely on recurrent neural networks (RNN). Particularly well suited to capture these long-range

temporal dependencies are models based on long short-term memory (LSTM) units [6] and variants thereof [7–9]. Thanks to automatic differentiation, parallel use of GPUs, and software packages like `theano` [10] or `torch` [11], it has become possible to easily fit such models to large scale data and obtain impressive results on tasks like text translation [12], speech recognition [13] and text or code generation [14, 15]. For algorithmic composition of monophonic music, it has been observed that RNN models based on LSTM units [16, 17] can capture long-term temporal dependencies far better than RNNs with simple units [18, 19].

Closest to our approach is the work of Eck and Schmidhuber [16] and Franklin [17]. Eck and Schmidhuber used an LSTM-RNN to generate a fixed Blues chord progression together with well-fitting improvised melodies. To model polyphonic music they discretized time into bins of equal duration, which has the disadvantage of using the same representation for one long note, e.g. a half note, and repeated notes of the same pitch, e.g. two quarter notes. Franklin experimented with different representations of pitch and duration and used an LSTM-RNN to reproduce a single long melody. These studies showed that LSTM-RNN’s are well suited to capture long-range temporal dependencies in music, but they did not demonstrate that these models can autonomously generate musically convincing novel melodies of a given style extracted from large datasets.

Here, we present a deep (multi-layer) model for algorithmic composition of monophonic melodies, based on RNN with gated recurrent units (GRU). We selected GRUs because they are simpler than LSTM units, but at least equally well suited to capture long-range temporal dependencies when used in RNNs (see [8, 20] for comparison studies). Our model represents monophonic music as a sequence of notes, where each note is given by its duration and pitch. Contrary to earlier studies, which evaluated LSTM based models of music only on small and artificial datasets [16, 17], we train our model on a large dataset of Irish folk songs. We apply the trained model to two different tasks: proposing convincing continuations of a melody and composing entire new songs.

2 Methods

To devise a statistical model that is able to complete and autonomously produce melodies, we train multi-layer RNNs on a large corpus of Irish folk songs. In the training phase, the network parameters are updated in order to accurately predict each upcoming note given the previously presented notes of songs in a training set (Fig. 2A). The model, once trained, can then be used to generate the upcoming notes itself, consequently producing whole new musical sequences.

After introducing our representation of melodies, we give a short introduction to recurrent neural networks and present the model and the training modalities. Finally, we explain how the model is used for algorithmic composition.

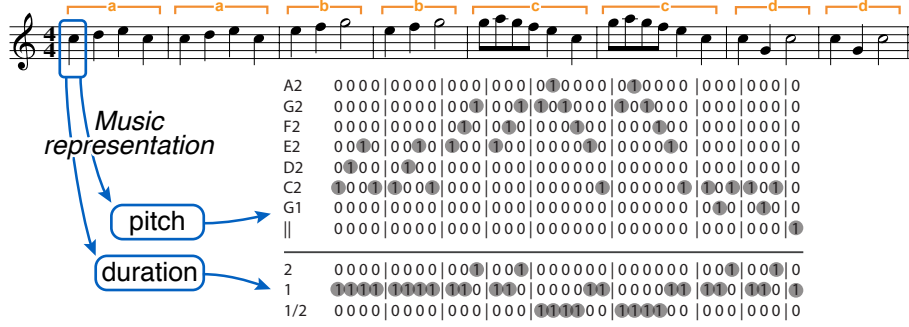


Fig. 1. Representation of a simple melody: The nursery rhyme *Brother John* in symbolic music notation, represented by pitch and duration matrices. Each column represents a single note. The four patterns composing the song are labeled by a, b, c and d.

2.1 Music Representation

Melodies are represented as sequences of notes, where each note is a combination of its pitch and duration. As such, each note n can be represented by two corresponding one-hot vectors (with only one entry of 1, and 0 for all others) of pitch $\mathbf{p}[n]$ and duration $\mathbf{d}[n]$. These vectors encode pitches and durations present in the training set. In addition, we include “song ending” and “silence” features as supplementary dimensions of the pitch vector. Any song can thus be mapped to a matrix of pitches and a second matrix of durations.

To reduce dictionary sizes as well as redundancies in representations, we chose to normalize both melody and duration. Melodies are normalized by transposing every song into C Major/A minor. Durations are normalized as relative to the most common duration in each song: if, for example, the most common duration in a song is the quarter note, we would represent eighth notes as “1/2”. For an example of a melody in this representation see Fig. 1.

2.2 A Brief Introduction to Recurrent Neural Networks

Artificial neural networks have a long history in machine learning, artificial intelligence and cognitive sciences (see [21] for a textbook, [22] for recent advances, [23] for an in-depth historical overview). Here we give a brief introduction for readers unfamiliar with the topic.

Artificial neural networks are non-linear functions $\mathbf{y} = f_{\mathbf{w}}(\mathbf{x})$, where input \mathbf{x} , output \mathbf{y} and the parameters (weights) \mathbf{w} can be elements of a high-dimensional space. A simple example of an artificial neural network with 2-dimensional in- and outputs is given by $y_1 = \tanh(w_{11}x_1 + w_{12}x_2)$ and $y_2 = \tanh(w_{21}x_1 + w_{22}x_2)$, which we write in short as $\mathbf{y} = \tanh(\mathbf{w}\mathbf{x})$. Characteristic of artificial neural networks is that the building blocks consist of a non-linear function σ (like \tanh) applied to a linear function $w_{11}x_1 + w_{12}x_2$, which is an abstraction of

the operation of biological neurons. If these building blocks are nested, e.g. $\mathbf{y} = \sigma_3(\mathbf{w}_3\sigma_2(\mathbf{w}_2\sigma_1(\mathbf{w}_1\mathbf{x})))$, one speaks of multi-layer (or deep) neural networks, with layer-specific weights $(\mathbf{w}_1, \mathbf{w}_2, \dots)$ and non-linearities $(\sigma_1, \sigma_2, \dots)$.

Deep neural networks are of particular interest for the approximation of high-dimensional and non-linear functions. For example, the function of recognizing object i in photos can be approximated by adjusting the weights such that output y_i is 1 if and only if an image \mathbf{x} of object i is given [22]. Formally, the weights can be adjusted to minimize a cost function, like the averaged square loss between target and output $\mathcal{L}(\mathbf{w}) = \frac{1}{S} \sum_{s=1}^S (\mathbf{y}_s - f_{\mathbf{w}}(\mathbf{x}_s))^2$ for some known input-output pairs $(\mathbf{x}_s, \mathbf{y}_s)$. Since artificial neural networks are differentiable in the parameters \mathbf{w} , this cost function is also differentiable and the parameters can be adjusted by changing them in direction of the gradient of the cost function $\Delta\mathbf{w} \propto \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w})$.

In recurrent neural networks (RNN) the inputs and outputs are sequences of arbitrary length and dimension. A simple example of a recurrent neural network with one hidden layer is given by $\mathbf{h}[n] = \sigma(\mathbf{w}_{xh}\mathbf{x}[n] + \mathbf{w}_{hh}\mathbf{h}[n-1])$ and $\mathbf{y}[n] = \sigma(\mathbf{w}_{hy}\mathbf{h}[n])$, where $\mathbf{x}[n]$, $\mathbf{h}[n]$, $\mathbf{y}[n]$ is the n -th element of the input, hidden, output sequence, respectively. This network is recurrent, since each hidden state $\mathbf{h}[n]$ depends on the previous hidden state $\mathbf{h}[n-1]$ and, therefore, on all previous input elements $\mathbf{x}[1], \mathbf{x}[2], \dots, \mathbf{x}[n]$. While these recurrent neural networks can in principle capture long-range temporal dependencies, they are difficult to fit to data by gradient descent, since the gradient involves the recurrent weights w_{hh} raised to high powers, which vanishes or explodes depending on the largest eigenvalue of w_{hh} [6]. This problem can be avoided by a reparametrization of the recurrent neural network (LSTM [6], GRU [9], other variants [20]). In Equation 1 we give the update equations for the GRU used in this study.

2.3 Model

To model distributions of pitch and duration, we use two separate multi-layer RNNs (see Fig. 2A) referred to as the *rhythm* and *melody* networks, respectively. The numbers of output units are equal to the dictionary sizes of our music representation. The three hidden layers of both the pitch and duration RNNs are composed of 128 GRUs [9] each, and are connected as shown in Fig. 2B. The model was implemented using the `theano` library [10].

For each note n in a musical sequence, the duration vector $\mathbf{d}[n]$ of this note is presented to the rhythm network, while the melody network receives both the pitch vector $\mathbf{p}[n]$ as well as the duration vector of the *upcoming* note $\mathbf{d}[n+1]$ as inputs. Each time a note is fed to the model, all internal states are updated and the rhythm network output gives a distribution over the possible upcoming durations $Pr(\mathbf{d}[n+1]|\mathbf{d}[n])$. In the same way, the melody network output gives a distribution over possible upcoming pitches $Pr(\mathbf{p}[n+1]|\mathbf{p}[n], \mathbf{d}[n+1])$.

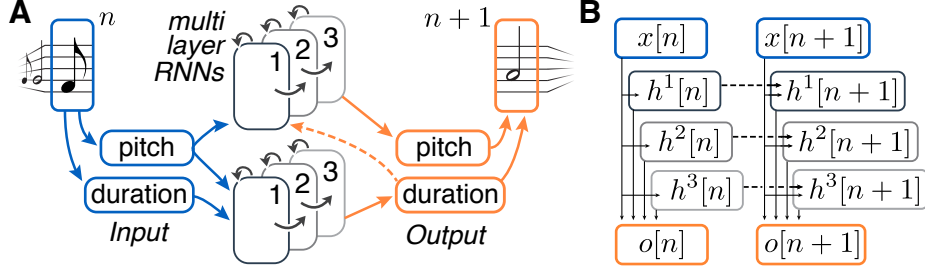


Fig. 2. Network architecture and connectivity: **A** Each note in a melody is separated into pitch and duration components, which are iteratively fed into two multilayer RNNs as inputs. Network outputs give probability distributions for pitch and duration of the next note. The rhythm network receives the current duration and pitch while the melody network receives the current pitch and the upcoming duration as input (dashed line). **B** The input layer x projects to all hidden layers h , as well as the output layer o . The hidden layers h project recurrently between notes (dashed lines) as well as feeding forward to all higher hidden layers and the output layer.

The update equations for the vector of layer activations $\mathbf{h}^i[n]$, update gates $\mathbf{z}^i[n]$ and reset gates $\mathbf{r}^i[n]$ at note n for layer $i \in \{1, 2, 3\}$ are given by

$$\mathbf{h}^i[n] = \mathbf{z}^i[n] \odot \mathbf{h}^i[n-1] + (1 - \mathbf{z}^i[n]) \odot \tilde{\mathbf{h}}^i[n], \quad (1)$$

$$\tilde{\mathbf{h}}^i[n] = \tanh \left(\mathbf{w}_{y^i h^i} \mathbf{y}^i[n] + \mathbf{r}^i[n] \odot \mathbf{w}_{h^i h^i} \mathbf{h}^i[n-1] \right), \quad (2)$$

$$\mathbf{z}^i[n] = \sigma \left(\mathbf{w}_{y^i z^i} \mathbf{y}^i[n] + \mathbf{w}_{h^i z^i} \mathbf{h}^i[n-1] + \mathbf{b}_{z^i}^i \right), \quad (3)$$

$$\mathbf{r}^i[n] = \sigma \left(\mathbf{w}_{y^i r^i} \mathbf{y}^i[n] + \mathbf{w}_{h^i r^i} \mathbf{h}^i[n-1] + \mathbf{b}_{r^i}^i \right), \quad (4)$$

where $\sigma(x) = (1 + \exp(x))^{-1}$ is the logistic sigmoid function, \odot denotes the element-wise product and \mathbf{y}^i is the feed-forward input to layer i , which consists of both the global inputs $\mathbf{x}[n]$ as well as hidden layer activations $\mathbf{h}^{j < i}[n]$ (see Fig. 2B). The update equation for the output unit i activation $\mathbf{o}_i[n]$ at note n is

$$\mathbf{o}_j[n] = \Theta(\mathbf{w}_{y^o o} \mathbf{y}^o[n] + \mathbf{b}^o)_j, \quad (5)$$

where \mathbf{y}^o is the feed-forward input of the output layer and $\Theta(\mathbf{x})_j = \frac{e^{x_j}}{\sum_k e^{x_k}}$ is the Softmax function. The Softmax normalization ensures that the values of the output units sum to one, which allows us to interpret the output of the two RNNs as probability distributions over pitches and over durations. For example, the probability of pitch j of an upcoming note is then given by

$$Pr(\mathbf{p}_j[n+1] = 1 \mid \text{previous notes and } \theta) = \mathbf{o}_j^{Melody}[n], \quad (6)$$

where the conditioning on previous notes and the model parameters θ highlights that the network output $\mathbf{o}[n]$ depends on them.

2.4 Training & Melody Generation

During training, the log-likelihood of model parameters θ for the rhythm and melody networks are separately optimized by stochastic gradient ascent with adaptive learning rate [24] ($\alpha = 10^{-3}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$). The log-likelihood of model parameters θ given the training songs is given by

$$\mathcal{L}(\theta) = \frac{1}{S} \sum_{s=1}^S \frac{1}{N_s - 1} \sum_{n=1}^{N_s-1} \log \left(Pr(\mathbf{x}_j^s[n+1] = 1 \mid \text{previous notes and } \theta) \right), \quad (7)$$

where S is the total number of songs, N_s the length of song s and $\mathbf{x}^s[n]$ is the duration vector $\mathbf{d}^s[n]$ for the rhythm network and the pitch vector $\mathbf{p}^s[n]$ for the melody network. The trainable model parameters θ are the connection matrices \mathbf{w}_{ab} for $a \in \{y^i, h^i\}$, $b \in \{h^i, z^i, r^i\}$ and $\mathbf{w}_{y^o o}$, the gate biases \mathbf{b}_z^i and \mathbf{b}_r^i , the output unit biases \mathbf{b}^o and the initial state of the hidden units $\mathbf{h}^i[0]$.

Networks are trained on 80% of songs of the musical corpus and tested on the remaining 20%. A training epoch consists of optimizing model parameters on each of randomly selected 200 melodies from the training set, where parameters are updated after each song. After each epoch, the model performances are evaluated on a random sample of 200 melodies from the testing and the training set. The parameters that minimize the likelihood on unseen data from the testing set are saved and used as final parameters.

Melody generation is achieved by closing the output-input loop of Eqs. 1-5. In each time step a duration and a pitch are sampled from the respective probability distributions and used as inputs in the next time step.

3 Results

3.1 Music Corpus

To automatically convert music corpora to our representation, we developed a toolbox that parses symbolic music files written in the abc notation¹ and converts them into our music representation. The results presented here are based on the Irish music corpus of Henrik Norbeck². It contains 2158 Irish tunes with an averaged length of 136 ± 84 notes. The basic statistics of pitch and duration transitions are shown in Fig. 3. Due to our normalization procedure, the most common duration is 1 and there is consequently a high probability of transition to the unitary duration from any state. Otherwise, shorter durations are mostly followed by the same or a complementary value, e.g. “3/2” followed by “1/2”. The most common pitches belong to the diatonic C Major scale. Transition from natural B, respectively G sharp, have high probability of ending in the closest C, respectively A, a property inherent to western music in order to resolve a melody.

¹ <http://abcnotation.com/>

² <http://www.norbeck.nu/abc/>

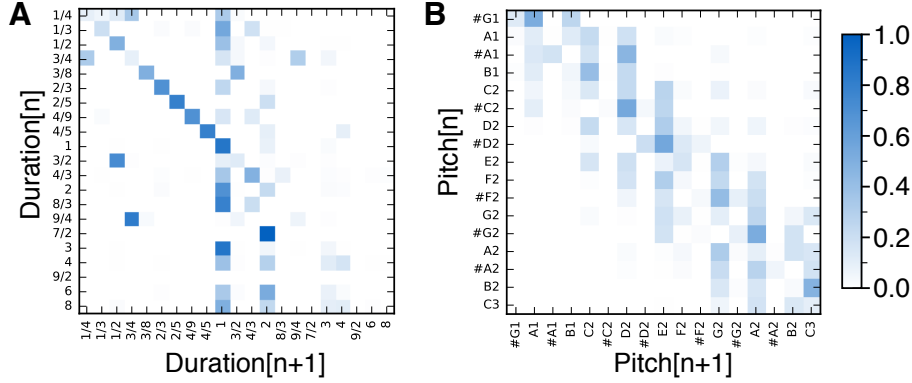


Fig. 3. Transition probabilities in the music corpus: **A** Probability of transitions from durations at note n to note $n + 1$ color graded from zero (white) to one (blue). **B** Same as in panel A, for transitions of a selected range of pitches.

The analysis of the relation between the pitch and duration features revealed that they are dependent, as expected from music theory. Therefore, we explicitly modeled the distribution over upcoming pitches as depending on the upcoming duration (dashed line in Fig. 2A), effectively splitting the joint distribution over note duration and pitch into conditional probabilities.

3.2 Song Continuation

To study song continuations, we present as input to the trained model the beginnings of previously unseen songs (the *seed*) and observe several continuations that our model produces (see Fig. 4). From a rhythmical point of view, it is interesting to notice that, even though the model had no notion of bars implemented, the metric structure was preserved in the generated continuations. Analyzing the continuations in Fig. 4A, we see that the A pentatonic scale that characterizes the seed is maintained everywhere except for continuation number 2. This is noteworthy, since the model is trained on a dataset that does not only contain music based on the pentatonic scale. Moreover, the rhythmic patterns of the seed are largely maintained in the continuations. Rhythmical patterns that are extraneous to the seed are also generated. For example, the pattern n_1 , which is the inversion of pattern a , can be observed several times. Importantly, the alternating structure of the seed (*abacabad*) is not present in the generated continuations. This indicates that the model is not able to capture this level of hierarchy. While less interesting than the first example from the rhythmical point of view, the seed presented in Fig. 4B, is clearly divided in a tonic area and a dominant area. Different continuations are coherent in the sense that they also alternate between these two areas, while avoiding the subdominant area almost everywhere.

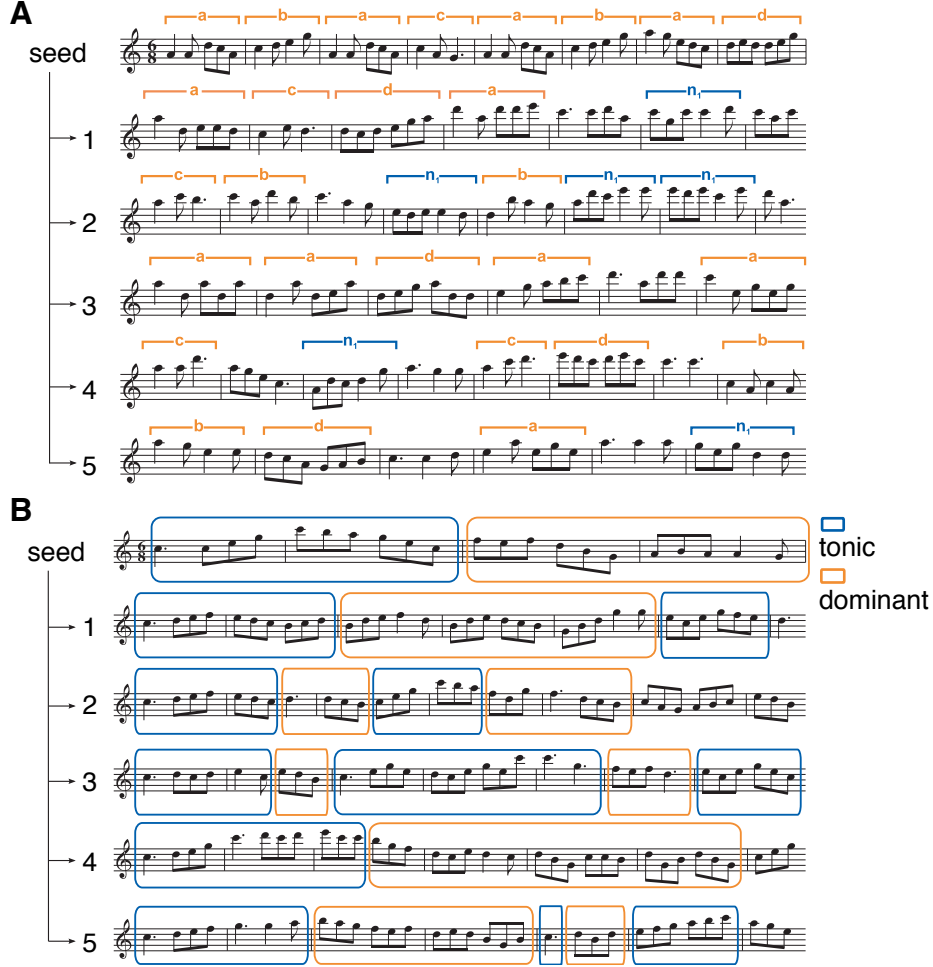


Fig. 4. Example of melody continuation and analysis: **A** The first line (seed) is presented as input to the model. The next five lines are five different possible continuations generated by our model. Rhythmical patterns present in the seed are labeled a, b, c and d. The label n_1 points at a new rhythmical pattern generated by the model. Unlabeled bars show other novel rhythmical patterns that appear only once. **B** A second example of song continuation, with analysis of tonal areas.

3.3 Autonomous Song Generation

Here, we assume that the model is able to learn the temporal dependencies in the musical examples over all timescales and use it to autonomously generate new pieces of music according to those dependencies. For the results presented here, we manually set the first two notes of the melody before generating notes by sampling from the output distributions. This step can be automatized but is



Fig. 5. Example of an autonomously generated Irish tune: See text for details. In this example a coherent temporary modulation to A minor was generated (bar 26). Worth noticing are passages that are reminiscent of the beginning at bar 13 and at bar 23.

needed in order to observe a distribution over possible upcoming notes at the model output. We iteratively add notes to the generated musical sequence until the “song ending” output is sampled.

We observe that the generated melodies (one example is shown in Fig. 5) are different from those found in the training set while carrying the same features on many time scales. Interestingly, the model is not explicitly aware of the time signature of each song, but bars can be easily added *a posteriori* indicating that the long-range rhythmical structure is perfectly learned and reproduced by our model. The generated melodies are produced according to production rules extracted from the training set, effectively creating an original composition that carries the features of the examples in the training dataset. The model is consequently able to generate new pieces of music in a completely autonomous manner.

4 Discussion

We fitted a statistical model that uses a recurrent neural network to a dataset of 2158 Irish melodies. Due to its recurrent connections and multiplicative units, the model is able to capture long-range temporal structure. The model contains no prior knowledge about a musical style but extracts all relevant features directly from the data. It can therefore be readily applied to other datasets of different musical styles. For example, training the model on the Nottingham Music Database³ yielded similar performance (data not shown).

We studied the model in two different settings. As a tool for composers, it can provide song continuations that are coherent with the beginning of the song both in terms of pitches and in terms of rhythmical patterns. The model also

³ <http://abc.sourceforge.net/NMD/>

allows the autonomous composition of new and complete musical sequences. The generated songs exhibit coherent metrical structure, in some cases temporary modulations to related keys, and are in general pleasant to hear.

Using RNNs for algorithmic composition allows to overcome the limitations of Markov chains in learning the long-range temporal dependencies of music [18]. A different class of models, namely artificial grammars, are naturally suited to generate these long-range dependencies due to their hierarchical structure. However, artificial grammars have proven to be much harder to learn from data than RNNs [25]. Therefore, researchers and composers usually define their own production rules in order to generate music [2]. Attempts have been made to infer context-free grammars [26] but applications to music are restricted to simple cases [27].

Evolutionary algorithms constitute another class of approaches that has been very popular in algorithmic composition [2]. They require the definition of a fitness function, i.e. a measure of the quality of musical compositions (the *individuals*, in the context of evolutionary algorithms). Based on the fitness function, the generative process corresponding to the best individuals is favoured and can undergo random mutations. This type of optimization process is not convenient for generative models for which gradient information is available, like neural networks. However, it can be used in rule-based generative models for which training is hard [28]. A common problem with evolutionary algorithms in music composition is that the definition of the fitness function is arbitrary and requires some kind of evaluation of the musical quality. This resulted in fitness functions often very specific to a certain style. However, similarly to neural networks, fitness functions can be defined based on the statistical similarity of the individuals to a target dataset of compositions [29].

Because of the ease of fitting these models to data as well as their expressiveness, learning algorithmic composition with recurrent neural networks seems promising. However, further quantitative evaluations are desirable. We see two quite different approaches for further evaluation. First, following standard practice in machine learning, our model could be compared to other approaches in terms of generalization, which is measured as the conditional probability of held-out test data (songs) of similar style. Second, the generated songs could be evaluated by humans in a Turing test setting or with a questionnaire following the SPECS methodology [30]. Another direction for interesting future work would be to fit a model to a corpus of polyphonic music and examining the influence of different representations.

Acknowledgments Research was supported by the Swiss National Science Foundation (200020_147200) and the European Research Council grant no. 268689 (MultiRules).

References

1. Ada Lovelace. Notes on L Menabrea’s “sketch of the analytical engine invented by Charles Babbage, esq.”. *Taylor’s Scientific Memoirs*, 3, 1843.

2. Jose D Fernández and Francisco Vico. Ai methods in algorithmic composition: A comprehensive survey. *Journal of Artificial Intelligence Research*, pages 513–582, 2013.
3. Kevin Jones. Compositional applications of stochastic processes. *Computer Music Journal*, 5(2):45–61, 1981.
4. Charles Ames. The markov process as a compositional model: a survey and tutorial. *Leonardo*, pages 175–187, 1989.
5. George Papadopoulos and Geraint Wiggins. Ai methods for algorithmic composition: A survey, a critical view and future prospects. In *AISB Symposium on Musical Creativity*, pages 110–117. Edinburgh, UK, 1999.
6. Sepp Hochreiter. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(2):107–116, 1998.
7. Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with lstm. *Neural computation*, 12(10):2451–2471, 2000.
8. Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
9. Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
10. James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a cpu and gpu math expression compiler. In *Proceedings of the Python for scientific computing conference (SciPy)*, volume 4, page 3. Austin, TX, 2010.
11. Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, number EPFL-CONF-192376, 2011.
12. Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
13. Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6645–6649. IEEE, 2013.
14. Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
15. Andrej Karpathy. The unreasonable effectiveness of recurrent neural networks. <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>, 2015. [Online; accessed 1-April-2016].
16. Douglas Eck and Jurgen Schmidhuber. Finding temporal structure in music: Blues improvisation with lstm recurrent networks. In *Neural Networks for Signal Processing, 2002. Proceedings of the 2002 12th IEEE Workshop on*, pages 747–756. IEEE, 2002.
17. Judy A Franklin. Computational models for learning pitch and duration using lstm recurrent neural networks. In *Proceedings of the Eighth International Conference on Music Perception and Cognition (ICMPC8)*, Adelaide, Australia. Causal Productions, 2004.
18. Peter M Todd. A connectionist approach to algorithmic composition. *Computer Music Journal*, 13(4):27–43, 1989.

19. Michael C Mozer. Neural network music composition by prediction: Exploring the benefits of psychoacoustic constraints and multi-scale processing. *Connection Science*, 6(2-3):247–280, 1994.
20. Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2342–2350, 2015.
21. Yoshua Bengio, Ian J Goodfellow, and Aaron Courville. Deep learning. *An MIT Press book in preparation. Draft chapters available at <http://www.iro.umontreal.ca/bengioy/dlbook>*, 2015.
22. Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
23. Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015.
24. Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
25. E Mark Gold. Language identification in the limit. *Information and control*, 10(5):447–474, 1967.
26. Craig G Nevill-Manning and Ian H Witten. Identifying hierarchical structure in sequences: A linear-time algorithm. *J. Artif. Intell. Res.(JAIR)*, 7:67–82, 1997.
27. Kris Makoto Kitani and Hideki Koike. Improvgenerator: Online grammatical induction for on-the-fly improvisation accompaniment. In *NIME*, pages 469–472, 2010.
28. Palle Dahlstedt. Autonomous evolution of complete piano pieces and performances. In *Proceedings of the ECAL Workshop on Music and Artificial Life, Lisbon, Portugal*, page 10, 2007.
29. Manuel Alfonseca, Manuel Cebrián, and Alfonso Ortega. Evolving computer-generated music by means of the normalized compression distance. *WSEAS Transactions on Information Science and Applications*, 2(9):343–348, 2005.
30. Anna Jordanous. A standardised procedure for evaluating creative systems: Computational creativity evaluation based on what it is to be creative. *Cognitive Computation*, 4(3):246–279, 2012.